

# 第三章

## 非球形扰动与广义最小二乘

# Outline

- 1 广义最小二乘
- 2 Feasible Generalized Least Squares
- 3 条件异方差检验
- 4 自相关检验
- 5 Heteroskedasticity and Autocorrelation Consistent (HAC) Covariance Matrix Estimator
- 6 Clustered Sampling

# 广义最小二乘

# Generalized Least Squares (GLS)

- Assumption 1.4 states that the  $n \times n$  matrix of conditional second moments  $\mathbb{E}(\varepsilon\varepsilon'|\mathbf{X})$  is spherical, that is, proportional to the identity matrix.
- If the error is not (conditionally) homoskedastic, the values of the diagonal elements of  $\mathbb{E}(\varepsilon\varepsilon'|\mathbf{X})$  are not the same.
- If there is correlation in the error term between observations, the values of the off-diagonal elements are not zero.
- The model that results when Assumption 1.4 is replaced by

$$\mathbb{E}(\varepsilon\varepsilon'|\mathbf{X}) = \sigma^2\mathbf{V}(\mathbf{X}), \quad \mathbf{V}(\mathbf{X}) \text{ nonsingular and known}, \quad (1)$$

is called the **generalized regression model**.

# Consequence of Relaxing Assumption 1.4

- ① The Gauss-Markov Theorem no longer holds for the OLS estimator:

$$\begin{aligned}
 \text{Var}(\mathbf{b}|\mathbf{X}) &= \text{Var}(\mathbf{b} - \boldsymbol{\beta}|\mathbf{X}) \\
 &= \text{Var}(\mathbf{A}\boldsymbol{\varepsilon}|\mathbf{X}), \quad \mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\
 &= \mathbf{A}\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X})\mathbf{A}' \\
 &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}(\mathbf{X})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}.
 \end{aligned}$$

- ② The  $t$ -ratio is not distributed as the  $t$  distribution. Thus, the  $t$ -test is no longer valid. The same comments apply to the  $F$ -test.
- ③ The OLS estimator is still unbiased (small sample) or consistent (large sample). For example,

$$\mathbb{E}(\mathbf{b} - \boldsymbol{\beta}|\mathbf{X}) = \mathbb{E}(\mathbf{A}\boldsymbol{\varepsilon}|\mathbf{X}) = \mathbf{A}\mathbb{E}(\boldsymbol{\varepsilon}|\mathbf{X}) = \mathbf{0}.$$

# Efficient Estimation with Known $\mathbf{V}$

Let the characteristic roots and vectors of  $\mathbf{V}$  be arranged in a diagonal matrix  $\mathbf{\Lambda}$  and an orthogonal matrix  $\mathbf{C}$ . Since  $\mathbf{V}$  is symmetric and positive definite, then

$$\mathbf{V} = \mathbf{C}\mathbf{\Lambda}\mathbf{C}'.$$

Therefore,

$$\mathbf{V}^{-1} = \mathbf{P}'\mathbf{P},$$

where  $\mathbf{P} \equiv \mathbf{C}\mathbf{\Lambda}^{-1/2}$ .

Now consider creating a new regression model by transforming  $(\mathbf{y}, \mathbf{X}, \boldsymbol{\varepsilon})$  by  $\mathbf{P}$  as

$$\tilde{\mathbf{y}} \equiv \mathbf{P}\mathbf{y}, \quad \tilde{\mathbf{X}} \equiv \mathbf{P}\mathbf{X}, \quad \tilde{\boldsymbol{\varepsilon}} \equiv \mathbf{P}\boldsymbol{\varepsilon}.$$

- Assumption 1.1 for  $(\mathbf{y}, \mathbf{X}, \varepsilon)$  implies that  $(\tilde{\mathbf{y}}, \tilde{\mathbf{X}}, \tilde{\varepsilon})$  too satisfies linearity:

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \tilde{\varepsilon}. \quad (2)$$

- Strict exogeneity is satisfied because

$$\mathbb{E}(\tilde{\varepsilon} \mid \tilde{\mathbf{X}}) = \mathbb{E}(\tilde{\varepsilon} \mid \mathbf{X}) = \mathbf{P}\mathbb{E}(\varepsilon \mid \mathbf{X}) = \mathbf{0}.$$

- Because  $\mathbf{V}$  is positive definite, the no-multicollinearity assumption is also satisfied.
- Assumption 1.4 is satisfied for the transformed model because

$$\begin{aligned} \mathbb{E}(\tilde{\varepsilon}\tilde{\varepsilon}' \mid \tilde{\mathbf{X}}) &= \mathbb{E}(\tilde{\varepsilon}\tilde{\varepsilon}' \mid \mathbf{X}) = \mathbf{P}\mathbb{E}(\varepsilon\varepsilon' \mid \mathbf{X})\mathbf{P}' \\ &= \sigma^2\mathbf{PVP}' = \sigma^2\mathbf{I}_n. \end{aligned}$$

- $\tilde{\varepsilon} \mid \tilde{\mathbf{X}}$  is normal if  $\varepsilon \mid \mathbf{X}$  is normal.

The Gauss-Markov Theorem for the transformed model implies that the BLUE of  $\beta$  for the generalized regression model is the OLS estimator applied to (2):

$$\hat{\beta}_{\text{GLS}} = (\tilde{\mathbf{X}}'\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}'\tilde{\mathbf{y}} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}.$$

This is the GLS estimator.

### Theorem (finite-sample properties of GLS)

- (a) (*Unbiasedness*) Under Assumptions 1.1-1.3,  $\mathbb{E}(\hat{\beta}_{\text{GLS}}|\mathbf{X}) = \beta$ .
- (b) (*Expression for the variance*) Under Assumptions 1.1-1.3 and the assumption (1),  $\text{Var}(\hat{\beta}_{\text{GLS}}|\mathbf{X}) = \sigma^2 \cdot (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}$ .
- (c) (*Efficiency of GLS*) Under the same set of assumptions as in (b), the GLS estimator is efficient in that the conditional variance of any unbiased estimator that is linear in  $\mathbf{y}$  is greater than or equal to  $\text{Var}(\hat{\beta}_{\text{GLS}}|\mathbf{X})$  in the matrix sense.



Proof.

(a) Note that  $\hat{\beta}_{\text{GLS}} - \beta = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\epsilon$ . Therefore,

$$\mathbb{E}(\hat{\beta}_{\text{GLS}} - \beta | \mathbf{X}) = \mathbb{E}[(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\epsilon | \mathbf{X}] = \mathbf{0}.$$

(b)

$$\begin{aligned} \text{Var}(\hat{\beta}_{\text{GLS}} | \mathbf{X}) &= \text{Var}(\hat{\beta}_{\text{GLS}} - \beta | \mathbf{X}) \\ &= (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\text{Var}(\epsilon | \mathbf{X})\mathbf{V}^{-1}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \\ &= \sigma^2 \cdot (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}. \end{aligned}$$

(c) Efficiency of GLS is implied by the fact that the transformed model satisfies Assumptions 1.1-1.4. □

## A Special Case: Weighted Least Squares (WLS)

When there is no correlation in the error term between observations,  $\mathbf{V}$  is diagonal. Let  $v_i(\mathbf{X})$  be the  $i$ -th diagonal element of  $\mathbf{V}(\mathbf{X})$ . So

$$\mathbb{E}(\varepsilon_i^2 \mid \mathbf{X}) = \sigma^2 \cdot v_i(\mathbf{X}).$$

It is easy to see that  $\mathbf{P}$  is also diagonal, with  $1/\sqrt{v_i(\mathbf{X})}$  in the  $i$ -th diagonal. Thus,  $(\tilde{\mathbf{y}}, \tilde{\mathbf{X}})$  is given by

$$\tilde{y}_i = \frac{y_i}{\sqrt{v_i(\mathbf{X})}}, \quad \tilde{\mathbf{x}}_i = \frac{\mathbf{x}_i}{\sqrt{v_i(\mathbf{X})}} \quad (i = 1, 2, \dots, n).$$

Therefore, efficient estimation under a known form of heteroskedasticity is first to weight each observation by the reciprocal of the square root of the variance  $v_i(\mathbf{X})$  and then apply OLS. This is called **the weighted regression (or the weighted least squares (WLS))**.

An important further special case is the case of a random sample where  $\{y_i, \mathbf{x}_i\}$  is i.i.d. across  $i$ . In this case,

- the error is unconditionally homoskedastic
- the error can be conditionally heteroskedastic
- $\mathbb{E}(\varepsilon_i^2 \mid \mathbf{X})$  depends only on  $\mathbf{x}_i$ , and the functional form of  $E(\varepsilon_i^2 \mid \mathbf{x}_i)$  is the same across  $i$ .

Thus

$$v_i(\mathbf{X}) = v(\mathbf{x}_i).$$

So the knowledge of  $\mathbf{V}(\cdot)$  comes down to a single function of  $K$  variables,  $v(\cdot)$ .

## Remark

- ① The finite-sample properties of GLS rest on the assumption that the regressors in the generalized regression model are strictly exogenous.
- ② If the regressors are not strictly exogenous but are merely predetermined, the GLS procedure to correct for serial correlation can make the estimator inconsistent.
- ③ If we do not know the function  $V(X)$ , we can estimate its functional form from the sample. This approach is called the **Feasible Generalized Least Squares (FGLS)**. Very little is known about the finite-sample properties of the FGLS estimator.
- ④ For FGLS, with an unrestricted matrix  $V(X)$ , there are  $n(n+1)/2$  additional parameters to estimate. This number is far too many to estimate with  $n$  observations. Some structure must be imposed on the model.

# Feasible Generalized Least Squares

To simplify the discussion, we assume that  $\{y_i, \mathbf{x}_i\}$  is i.i.d., and

$$\mathbb{E}(\varepsilon_i^2 \mid \mathbf{x}_i) = \mathbf{z}_i' \boldsymbol{\alpha}, \quad (3)$$

where  $\mathbf{z}_i$  is a function of  $\mathbf{x}_i$ .

**WLS with known  $\boldsymbol{\alpha}$ :** The WLS estimator is given by

$$\hat{\boldsymbol{\beta}}(\mathbf{V}) \equiv \left( \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{y}_i = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y},$$

with

$$\mathbf{V} \equiv \begin{bmatrix} \mathbf{z}_1' \boldsymbol{\alpha} & & \\ & \ddots & \\ & & \mathbf{z}_n' \boldsymbol{\alpha} \end{bmatrix}.$$

If Assumption 2.3 is strengthened by the condition that

$$\mathbb{E}(\varepsilon_i \mid \mathbf{x}_i) = 0,$$

then  $\mathbb{E}(\tilde{\varepsilon}_i \mid \tilde{\mathbf{x}}_i) = 0$ . Therefore, provided that  $\mathbb{E}(\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i')$  is nonsingular, the WLS estimator is consistent and asymptotically normal (Assumptions 2.1-2.5 are satisfied), and the asymptotic variance is

$$\begin{aligned} \text{Avar}(\hat{\beta}(\mathbf{V})) &= \mathbb{E}(\tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i')^{-1} \\ &= \text{plim} \left( \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i' \right)^{-1} \\ &= \text{plim} \left( \frac{1}{n} \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \right)^{-1}. \end{aligned}$$

Regression of  $e_i^2$  on  $\mathbf{z}_i$  provides a consistent estimate of  $\alpha$ : If  $\alpha$  is unknown, it can be estimated by running a separate regression.

- Note that (3) can be written as a regression equation:

$$\varepsilon_i^2 = \mathbf{z}_i' \alpha + \eta_i, \quad (4)$$

where  $\eta_i \equiv \varepsilon_i^2 - \mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i)$ . By construction,  $\mathbb{E}(\eta_i | \mathbf{x}_i) = 0$ , which implies that  $\mathbb{E}(\eta_i | \mathbf{z}_i) = 0$ .

- Hence, provided that  $\mathbb{E}(\mathbf{z}_i \mathbf{z}_i')$  is nonsingular, the OLS estimator of  $\alpha$  is consistent and asymptotically normal.
- When  $\varepsilon_i$  is replaced by  $e_i$  in the regression (4), the OLS estimator, call it  $\hat{\alpha}$ , is consistent for  $\alpha$ .



WLS with Estimated  $\alpha$ 

**Step 1:** Estimate the equation  $y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i$  by OLS and compute the OLS residuals  $e_i$ .

**Step 2:** Regress  $e_i^2$  on  $\mathbf{z}_i$ , to obtain the OLS coefficient estimate  $\hat{\alpha}$ .

**Step 3:** Re-estimate the equation  $y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i$  by WLS, using  $1/\sqrt{\mathbf{z}'_i \hat{\alpha}}$  as the weight for observation  $i$ .

It can be shown that the WLS estimator is asymptotically more efficient than the OLS estimator.

The superiority of WLS over OLS, however, rests on the premise that the sample size is sufficiently large and the functional form of the conditional second moment is correctly specified.

# 条件异方差检验

# White Test

Recall that  $\hat{\mathbf{S}} = \frac{1}{n} \sum_{i=1}^n e_i^2 \mathbf{x}_i \mathbf{x}_i'$  is consistent for  $\mathbf{S} = \mathbb{E}(\varepsilon_i^2 \mathbf{x}_i \mathbf{x}_i')$ . Under conditional homoskedasticity,  $s^2 \mathbf{S}_{xx}$  is consistent for  $\mathbf{S}$ . So the difference between the two consistent estimators should vanish:

$$\hat{\mathbf{S}} - s^2 \mathbf{S}_{xx} = \frac{1}{n} \sum_{i=1}^n (e_i^2 - s^2) \mathbf{x}_i \mathbf{x}_i' \xrightarrow{p} \mathbf{0}.$$

Let  $\boldsymbol{\psi}_i$  be a vector collecting unique and nonconstant elements of the  $K \times K$  symmetric matrix  $\mathbf{x}_i \mathbf{x}_i'$ . Then,

$$\mathbf{c}_n \equiv \frac{1}{n} \sum_{i=1}^n (e_i^2 - s^2) \boldsymbol{\psi}_i \xrightarrow{p} \mathbf{0}.$$

Under some conditions appropriate for a CLT to be applicable,  $\sqrt{n} \mathbf{c}_n$  converges in distribution to a normal distribution with mean zero and some asymptotic variance  $\mathbf{B}$ , so for any consistent estimator  $\hat{\mathbf{B}}$  of  $\mathbf{B}$ ,

$$n \cdot \mathbf{c}_n' \hat{\mathbf{B}}^{-1} \mathbf{c}_n \xrightarrow{d} \chi^2(m),$$

where  $m$  is the dimension of  $\mathbf{c}_n$ .

For a certain choice of  $\hat{\mathbf{B}}$ , this statistic can be computed as  $nR^2$  from the following auxiliary regression:

$$\text{regress } e_i^2 \text{ on a constant and } \psi_i. \quad (5)$$

### Theorem (White's Test for Conditional Heteroskedasticity)

*In addition to Assumptions 2.1 and 2.4, suppose that (a)  $\{y_i, \mathbf{x}_i\}$  is i.i.d. with finite  $\mathbb{E}(\varepsilon_i^2 \mathbf{x}_i \mathbf{x}_i')$ , (b)  $\varepsilon_i$  is independent of  $\mathbf{x}_i$ , and (c) a certain condition holds on the moments of  $\varepsilon_i$  and  $\mathbf{x}_i$ . Then,*

$$nR^2 \xrightarrow{d} \chi^2(m),$$

*where  $R^2$  is the  $R^2$  from the auxiliary regression (5), and  $m$  is the dimension of  $\psi_i$ .*

# Breusch-Pagan Test

Consider the linear model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i, \quad (6)$$

where  $\varepsilon_i$  is **normally** and independently distributed with mean zero and variance

$$\sigma_i^2 = \sigma^2 \cdot h(\mathbf{z}_i' \boldsymbol{\alpha}),$$

where  $\boldsymbol{\alpha}$  is  $p \times 1$  vector and the first element is 1. The null hypothesis of homoscedasticity can be written as

$$\mathbb{H}_0 : \alpha_2 = \cdots = \alpha_p = 0.$$

It is also assumed that  $\mathbf{x}_i$  and  $\mathbf{z}_i$  are exogenous.

## Breusch-Pagan Test

**Step 1:** Run regression (6), obtain residuals  $e_i$ 's.

**Step 2:** Obtain  $g_i$  which is given by

$$g_i \equiv \frac{e_i^2}{s^2}, \quad \text{with } s^2 = \frac{\mathbf{e}'\mathbf{e}}{n}$$

for  $i = 1, \dots, n$ .

**Step 3:** Regress  $g_i$  on  $\mathbf{z}_i$ , obtain the explained sum of squares which is  $\mathbf{g}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{g} - n$ , where  $\mathbf{g} = (g_1, \dots, g_n)'$ ,  $\mathbf{Z} = (\mathbf{z}'_1, \dots, \mathbf{z}'_n)'$ .

**Step 4:** Construct the test statistic, under  $\mathbb{H}_0$ ,

$$\text{LM} \equiv \frac{1}{2} [\mathbf{g}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{g} - n] \xrightarrow{d} \chi^2(p-1).$$

# 自相关检验

- When the regressors include a constant (true in virtually all known applications), Assumption 2.5 implies that the error term is a scalar martingale difference sequence.
- If the error is found to be serially correlated, that is an indication of a failure of Assumption 2.5.
- Serial correlation has traditionally been an important subject in econometrics, and there are available a number of tests for serial correlation.
- However, some of them (e.g., Durbin-Watson test) require that the regressors be strictly exogenous.



## Digression: Box-Pierce and Ljung-Box

Suppose we have a sample  $\{z_1, \dots, z_n\}$  drawn from a scalar covariance-stationary process. The sample  $j$ -th order autocovariance is

$$\hat{\gamma}_j \equiv \frac{1}{n} \sum_{t=j+1}^n (z_t - \bar{z}_n)(z_{t-j} - \bar{z}_n),$$

where

$$\bar{z}_n \equiv \frac{1}{n} \sum_{t=1}^n z_t.$$

The sample  $j$ -th order autocorrelation coefficient,  $\hat{\rho}_j$ , is defined as

$$\hat{\rho}_j \equiv \frac{\hat{\gamma}_j}{\hat{\gamma}_0}.$$

If  $\{z_t\}$  is ergodic stationary, then it is easy to show that  $\hat{\gamma}_j$  is consistent for  $\gamma_j$ , and  $\hat{\rho}_j$  is consistent for  $\rho_j$ .

## Theorem

Suppose  $\{z_t\}$  can be written as  $\mu + \varepsilon_t$ , where  $\varepsilon_t$  is a stationary martingale difference sequence with “own” conditional homoskedasticity:

$$(\text{own conditional homoskedasticity}) \quad \mathbb{E}(\varepsilon_t^2 \mid \varepsilon_{t-1}, \varepsilon_{t-2}, \dots) = \sigma^2 > 0.$$

Then,

$$\sqrt{n}\hat{\gamma} \xrightarrow{d} N(\mathbf{0}, \sigma^4 \mathbf{I}_p) \quad \text{and} \quad \sqrt{n}\hat{\rho} \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p),$$

where  $\hat{\gamma} = (\hat{\gamma}_1, \dots, \hat{\gamma}_p)'$  and  $\hat{\rho} = (\hat{\rho}_1, \dots, \hat{\rho}_p)'$ . (Cf. [Hayashi, 2000], p.171, Q9)

For testing  $\mathbb{H}_0 : \rho_1 = \dots = \rho_p = 0$ , we have:

- Box-Pierce Q-test:

$$Q \equiv n \sum_{j=1}^p \hat{\rho}_j^2 = \sum_{j=1}^p (\sqrt{n}\hat{\rho}_j)^2 \xrightarrow{d} \chi^2(p).$$

- Ljung-Box Q-test:

$$Q^* \equiv n(n+2) \sum_{j=1}^p \frac{\hat{\rho}_j^2}{n-j} = \sum_{j=1}^p \frac{n+2}{n-j} (\sqrt{n}\hat{\rho}_j)^2 \xrightarrow{d} \chi^2(p).$$

# Sample Autocorrelations Calculated from Residuals

For the regression model described by Assumptions 2.1-2.5, if  $\varepsilon_t$  were observable, we would calculate

$$\tilde{\rho}_j \equiv \frac{\tilde{\gamma}_j}{\tilde{\gamma}_0} \quad (j = 1, 2, \dots),$$

where

$$\tilde{\gamma}_j \equiv \frac{1}{n} \sum_{t=j+1}^n \varepsilon_t \varepsilon_{t-j} \quad (j = 0, 1, 2, \dots).$$

Consider the realistic case where we do not observe  $\varepsilon_t$ . We can replace  $\varepsilon_t$  by  $e_t$  and calculate

$$\hat{\rho}_j \equiv \frac{\hat{\gamma}_j}{\hat{\gamma}_0} \quad (j = 1, 2, \dots),$$

where

$$\hat{\gamma}_j \equiv \frac{1}{n} \sum_{t=j+1}^n e_t e_{t-j} \quad (j = 0, 1, 2, \dots).$$

Is it all right to use residual-based  $Q$  test derived from  $\{\hat{\rho}_j\}$  for testing serial correlation?

Answer: YES! But only if the regressors are strictly exogenous.

To see this, recall that

$$e_t = \varepsilon_t - \mathbf{x}_t'(\mathbf{b} - \boldsymbol{\beta}).$$

Therefore,

$$\begin{aligned}\hat{\gamma}_j &\equiv \frac{1}{n} \sum_{t=j+1}^n e_t e_{t-j} = \frac{1}{n} \sum_{t=j+1}^n [\varepsilon_t - \mathbf{x}_t'(\mathbf{b} - \boldsymbol{\beta})][\varepsilon_{t-j} - \mathbf{x}_{t-j}'(\mathbf{b} - \boldsymbol{\beta})] \\ &= \tilde{\gamma}_j - \frac{1}{n} \sum_{t=j+1}^n (\mathbf{x}_{t-j} \varepsilon_t + \mathbf{x}_t \varepsilon_{t-j})'(\mathbf{b} - \boldsymbol{\beta}) \\ &\quad + (\mathbf{b} - \boldsymbol{\beta})' \left( \frac{1}{n} \sum_{t=j+1}^n \mathbf{x}_t \mathbf{x}_{t-j}' \right) (\mathbf{b} - \boldsymbol{\beta}).\end{aligned}\tag{7}$$

Since  $\mathbf{b} - \boldsymbol{\beta} \xrightarrow{p} \mathbf{0}$ ,  $\hat{\gamma}_j - \tilde{\gamma}_j \xrightarrow{p} 0$  ( $j = 0, 1, 2, \dots$ ), and the difference between  $\tilde{\rho}_j$  and  $\hat{\rho}_j$  vanishes in large samples.

However, by multiplying both sides of (7) by  $\sqrt{n}$ , we have

$$\begin{aligned}\sqrt{n}\hat{\gamma}_j &= \sqrt{n}\tilde{\gamma}_j - \frac{1}{n} \sum_{t=j+1}^n (\mathbf{x}_{t-j}\varepsilon_t + \mathbf{x}_t\varepsilon_{t-j})' \sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) \\ &\quad + \underbrace{\sqrt{n}(\mathbf{b} - \boldsymbol{\beta})' \left( \frac{1}{n} \sum_{t=j+1}^n \mathbf{x}_t \mathbf{x}_{t-j}' \right) (\mathbf{b} - \boldsymbol{\beta})}_{\xrightarrow{p} 0}.\end{aligned}$$

Regarding the second term, we have

$$\frac{1}{n} \sum_{t=j+1}^n (\mathbf{x}_{t-j}\varepsilon_t + \mathbf{x}_t\varepsilon_{t-j}) \xrightarrow{p} \mathbb{E}(\mathbf{x}_{t-j}\varepsilon_t) + \mathbb{E}(\mathbf{x}_t\varepsilon_{t-j}).$$

If the regressors are strictly exogenous in the sense that  $\mathbb{E}(\mathbf{x}_t\varepsilon_s) = \mathbf{0}$  for all  $t, s$ , then

$$\mathbb{E}(\mathbf{x}_{t-j}\varepsilon_t) + \mathbb{E}(\mathbf{x}_t\varepsilon_{t-j}) = \mathbf{0}.$$

So the difference between  $\sqrt{n}\tilde{\rho}_j$  and  $\sqrt{n}\hat{\rho}_j$  vanishes, which means that the  $Q$  statistic calculated from  $\{e_t\}$  is asymptotically chi-squared, and we can use this residual-based  $Q$  to test for serial correlation.

# Testing with Predetermined, but Not Strictly Exogenous, Regressors

When the regressors are not strictly exogenous, we need to modify the  $Q$  statistic to restore its asymptotic distribution. Consider two restrictions:

Stronger form of predeterminedness:

$$\mathbb{E}(\varepsilon_t \mid \varepsilon_{t-1}, \varepsilon_{t-2}, \dots, \mathbf{x}_t, \mathbf{x}_{t-1}, \dots) = 0. \quad (8)$$

Stronger form of conditional homoskedasticity:

$$\mathbb{E}(\varepsilon_t^2 \mid \varepsilon_{t-1}, \varepsilon_{t-2}, \dots, \mathbf{x}_t, \mathbf{x}_{t-1}, \dots) = \sigma^2 > 0. \quad (9)$$

Remark:

- (8) is stronger than Assumption 2.3 and implies that  $\mathbf{g}_t$  is an m.d.s.
- (9) is obviously stronger than Assumption 2.7.

## Theorem (testing for serial correlation with predetermined regressors)

Suppose that Assumptions 2.1, 2.2, 2.4, (8), and (9) are satisfied. Then,

$$\sqrt{n}\hat{\gamma} \xrightarrow{d} N(\mathbf{0}, \sigma^4 \cdot (\mathbf{I}_p - \Phi)) \quad \text{and} \quad \sqrt{n}\hat{\rho} \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p - \Phi),$$

where  $\hat{\gamma} = (\hat{\gamma}_1, \dots, \hat{\gamma}_p)'$ ,  $\hat{\rho} = (\hat{\rho}_1, \dots, \hat{\rho}_p)'$ , and  $\phi_{jk}$ , the  $(j, k)$ -th element of the  $p \times p$  matrix  $\Phi$ , is given by

$$\phi_{jk} = \mathbb{E}(\mathbf{x}_t \varepsilon_{t-j})' \mathbb{E}(\mathbf{x}_t \mathbf{x}_t')^{-1} \mathbb{E}(\mathbf{x}_t \varepsilon_{t-k}) / \sigma^2.$$

By the Ergodic Theorem, matrix  $\Phi$  is consistently estimated by

$$\hat{\Phi} \equiv (\hat{\phi}_{jk}), \quad \hat{\phi}_{jk} \equiv \bar{\boldsymbol{\mu}}_j' \mathbf{S}_{\mathbf{xx}}^{-1} \bar{\boldsymbol{\mu}}_k / s^2 \quad (j, k = 1, 2, \dots, p),$$

where

$$s^2 \equiv \frac{1}{n-K} \sum_{t=1}^n e_t^2, \quad \bar{\boldsymbol{\mu}}_j \equiv \frac{1}{n} \sum_{t=j+1}^n \mathbf{x}_t \cdot e_{t-j}.$$

It follows that

$$\text{modified Box-Pierce } Q \equiv n \cdot \hat{\rho}' (\mathbf{I}_p - \hat{\Phi})^{-1} \hat{\rho} \xrightarrow{d} \chi^2(p).$$

## Proof.

Note that

$$\begin{aligned}
\sqrt{n}\hat{\gamma}_j &= \sqrt{n}\tilde{\gamma}_j - \frac{1}{n} \sum_{t=j+1}^n (\mathbf{x}_{t-j}\varepsilon_t + \mathbf{x}_t\varepsilon_{t-j})' \sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) \\
&\quad + \sqrt{n}(\mathbf{b} - \boldsymbol{\beta})' \left( \frac{1}{n} \sum_{t=j+1}^n \mathbf{x}_t \mathbf{x}_{t-j}' \right) (\mathbf{b} - \boldsymbol{\beta}) \\
&\stackrel{a}{\sim} \sqrt{n}\tilde{\gamma}_j - \frac{1}{n} \sum_{t=j+1}^n (\mathbf{x}_{t-j}\varepsilon_t + \mathbf{x}_t\varepsilon_{t-j})' \sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) \\
&\stackrel{a}{\sim} \sqrt{n}\tilde{\gamma}_j - \mathbb{E}(\mathbf{x}_{t-j}\varepsilon_t + \mathbf{x}_t\varepsilon_{t-j})' \sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) \\
&= \sqrt{n}\tilde{\gamma}_j - \boldsymbol{\mu}_j' \sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) \quad \text{where } \boldsymbol{\mu}_j \equiv \mathbb{E}(\mathbf{x}_t\varepsilon_{t-j}) \\
&= \sqrt{n}\tilde{\gamma}_j - \boldsymbol{\mu}_j' \sqrt{n} \mathbf{S}_{\mathbf{xx}}^{-1} \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t \varepsilon_t \\
&= \sqrt{n} \frac{1}{n} \sum_{t=j+1}^n \varepsilon_t \varepsilon_{t-j} - \boldsymbol{\mu}_j' \sqrt{n} \mathbf{S}_{\mathbf{xx}}^{-1} \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t \varepsilon_t \\
&\stackrel{a}{\sim} \sqrt{n} \frac{1}{n} \sum_{t=1}^n \varepsilon_t \varepsilon_{t-j} - \boldsymbol{\mu}_j' \sqrt{n} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1} \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t \varepsilon_t = \mathbf{c}_j' \sqrt{n} \bar{\mathbf{g}}_j,
\end{aligned}$$



## Proof.

where

$$\mathbf{c}_j \equiv \begin{bmatrix} 1 \\ -\boldsymbol{\Sigma}_{\mathbf{x}\mathbf{x}}^{-1} \boldsymbol{\mu}_j \end{bmatrix}, \quad \bar{\mathbf{g}}_j \equiv \frac{1}{n} \sum_{t=1}^n \mathbf{g}_{jt}, \quad \mathbf{g}_{jt} \equiv \begin{bmatrix} \varepsilon_{t-j} \varepsilon_t \\ \mathbf{x}_t \varepsilon_t \end{bmatrix}.$$

Thus, letting  $\hat{\boldsymbol{\gamma}} = (\hat{\gamma}_1, \dots, \hat{\gamma}_p)'$ , we have proved that

$$\sqrt{n} \hat{\boldsymbol{\gamma}} \underset{a}{\sim} \mathbf{C}' \sqrt{n} \bar{\mathbf{g}}, \quad (10)$$

where

$$\mathbf{C} \equiv \begin{bmatrix} \mathbf{c}_1 & & \\ & \ddots & \\ & & \mathbf{c}_p \end{bmatrix}, \quad \bar{\mathbf{g}} \equiv \frac{1}{n} \sum_{t=1}^n \mathbf{g}_t, \quad \mathbf{g}_t \equiv \begin{bmatrix} \mathbf{g}_{1t} \\ \vdots \\ \mathbf{g}_{pt} \end{bmatrix}.$$

Using LIE and (8), we can show that  $\mathbf{g}_t$  is a m.d.s. Clearly,  $\mathbf{g}_t$  is also ergodic stationary. Therefore,

$$\sqrt{n} \bar{\mathbf{g}} \xrightarrow{d} N(\mathbf{0}, \mathbb{E}(\mathbf{g}_t \mathbf{g}_t')). \quad (11)$$

The next step is to calculate  $\mathbb{E}(\mathbf{g}_t \mathbf{g}_t')$ . Its  $(j, k)$  block is

$$\mathbb{E}(\mathbf{g}_{jt} \mathbf{g}_{kt}') = \begin{bmatrix} \mathbb{E}(\varepsilon_t^2 \varepsilon_{t-j} \varepsilon_{t-k}) & \mathbb{E}(\mathbf{x}_t' \varepsilon_t^2 \varepsilon_{t-j}) \\ \mathbb{E}(\mathbf{x}_t \varepsilon_t^2 \varepsilon_{t-k}) & \mathbb{E}(\varepsilon_t^2 \mathbf{x}_t \mathbf{x}_t') \end{bmatrix}.$$

## Proof.

By using LIE, (8), and (9), it is easy to show that

$$\mathbb{E}(\mathbf{g}_{jt}\mathbf{g}'_{kt}) = \begin{bmatrix} \sigma^4 \delta_{jk} & \sigma^2 \boldsymbol{\mu}'_j \\ \sigma^2 \boldsymbol{\mu}_k & \sigma^2 \boldsymbol{\Sigma}_{xx} \end{bmatrix},$$

where  $\delta_{jk}$  is 1 if  $j = k$  and 0 otherwise. From (10) and (11), we have

$$\text{Avar}(\hat{\gamma}) = \mathbf{C}' \mathbb{E}(\mathbf{g}_t \mathbf{g}'_t) \mathbf{C}.$$

Its  $(j, k)$  element is

$$\mathbf{c}'_j \mathbb{E}(\mathbf{g}_{jt} \mathbf{g}'_{kt}) \mathbf{c}_k = \sigma^4 \cdot [\delta_{jk} - \boldsymbol{\mu}'_j \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\mu}_k / \sigma^2].$$

So

$$\sqrt{n} \hat{\gamma} \xrightarrow{d} N(\mathbf{0}, \sigma^4 \cdot (\mathbf{I}_p - \boldsymbol{\Phi})).$$

Since  $\sqrt{n} \hat{\boldsymbol{\rho}} \underset{a}{\sim} \sqrt{n} \hat{\gamma} / \sigma^2$ , the limiting distribution of  $\sqrt{n} \hat{\boldsymbol{\rho}}$  is the same as that of  $\sqrt{n} \hat{\gamma} / \sigma^2$ , which is  $N(\mathbf{0}, \mathbf{I}_p - \boldsymbol{\Phi})$ . □

# An Auxiliary Regression-Based Test

Consider the following auxiliary regression:

$$\text{regress } e_t \text{ on } \mathbf{x}_t, e_{t-1}, e_{t-2}, \dots, e_{t-p}. \quad (12)$$

From this regression, we can calculate the  $F$  statistic for the hypothesis that the  $p$  coefficients of  $e_{t-1}, \dots, e_{t-p}$  are all zero.

- Under the hypothesis of the previous theorem, the modified  $Q$  statistic is asymptotically equivalent to  $p \cdot F$ , so  $p \cdot F$ , too, is asymptotically chi-squared. (Cf. [Hayashi, 2000], p.173, Q11)
- This  $p \cdot F$  statistic, in turn, is asymptotically equivalent to  $nR^2$  from the auxiliary regression. The test based on  $nR^2$  is called the Breusch-Godfrey test for serial correlation.

We now prove the last statement.

## Proof.

Recall that

$$p \cdot F = \frac{SSR_R - SSR_U}{SSR_U / (n - K - p)} = (n - K - p) \frac{SSR_R - SSR_U}{SSR_U},$$

where  $SSR_U$  is from (12) and  $SSR_R$  is from the following restricted regression:

$$\text{regress } e_t \text{ on } \mathbf{x}_t. \quad (13)$$

However, since  $e_t$  is the OLS residual,  $\mathbf{x}_t$  in (13) has no explanatory power. So  $SSR_R = e'e$ , and  $p \cdot F$  is numerically identical to

$$(n - K - p) \frac{R_{uc}^2}{1 - R_{uc}^2},$$

where  $R_{uc}^2$  is the uncentered  $R^2$  for (12).

## Proof.

Since the sample mean of  $e_t$  is zero, we also have

$$p \cdot F = (n - K - p) \frac{R^2}{1 - R^2},$$

where  $R^2$  is equal to the  $R^2$  for (12). Solving this equation for  $R^2$  and multiplying both sides by  $n$ , we obtain

$$nR^2 = \frac{n}{n - K - p} \cdot \frac{1}{1 + \frac{p \cdot F}{n - K - p}} \cdot p \cdot F.$$

It is easy to see that  $p \cdot F - nR^2 \xrightarrow[p]{} 0$ . Therefore,  $nR^2$  from (12) is asymptotically  $\chi^2(p)$ . □

# Heteroskedasticity and Autocorrelation Consistent (HAC) Covariance Matrix Estimator

# Asymptotics for Sample Means of Serially Correlated Processes

Consider the sample mean  $\bar{y} \equiv \frac{1}{n} \sum_{t=1}^n y_t$  for serially correlated process  $\{y_t\}$ .

Theorem (LLN for covariance-stationary processes with vanishing autocovariances)

Let  $\{y_t\}$  be covariance-stationary with mean  $\mu$  and  $\{\gamma_j\}$  be the autocovariances of  $\{y_t\}$ . Then

- (a)  $\bar{y} \xrightarrow{m.s.} \mu$  as  $n \rightarrow \infty$ , if  $\lim_{j \rightarrow \infty} \gamma_j = 0$ .
- (b)  $\lim_{n \rightarrow \infty} \text{Var}(\sqrt{n}\bar{y}) = \sum_{j=-\infty}^{\infty} \gamma_j < \infty$ , if  $\{\gamma_j\}$  is summable.

For a covariance-stationary process  $\{y_t\}$ , define the **long-run variance** to be the limit as  $n \rightarrow \infty$  of  $\text{Var}(\sqrt{n}\bar{y})$ .

Note that

$$\begin{aligned}
 \text{Var}(\sqrt{n}\bar{y}) &= \frac{1}{n} \text{Var} \left( \sum_{t=1}^n y_t \right) = \frac{1}{n} \left( n\gamma_0 + \sum_{t=1}^n \sum_{s \neq t}^n \text{Cov}(y_t, y_s) \right) \\
 &= \gamma_0 + \frac{2}{n} \sum_{t=1}^n \sum_{s < t}^n \text{Cov}(y_t, y_s) = \gamma_0 + \frac{2}{n} \sum_{j=1}^{n-1} (n-j) \gamma_j \\
 &= \gamma_0 + 2 \sum_{j=1}^{n-1} \left(1 - \frac{j}{n}\right) \gamma_j.
 \end{aligned}$$

### Theorem (CLT for MA( $\infty$ ))

Let  $y_t = \mu + \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}$  where  $\{\varepsilon_t\}$  is independent white noise and  $\sum_{j=0}^{\infty} |\psi_j| < \infty$ . Then

$$\sqrt{n}(\bar{y} - \mu) \xrightarrow{d} N \left( 0, \sum_{j=-\infty}^{\infty} \gamma_j \right).$$



## Theorem

Let  $\bar{\mathbf{y}} \equiv \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t$  be the sample mean of a vector process  $\{\mathbf{y}_t\}$ .

- (a)  $\bar{\mathbf{y}} \xrightarrow{m.s.} \boldsymbol{\mu}$  if each diagonal element of  $\boldsymbol{\Gamma}_j$  goes to zero as  $j \rightarrow \infty$ .
- (b)  $\lim_{n \rightarrow \infty} \text{Var}(\sqrt{n}\bar{\mathbf{y}}) = \sum_{j=-\infty}^{\infty} \boldsymbol{\Gamma}_j$ , if  $\{\boldsymbol{\Gamma}_j\}$  is summable.
- (c) If  $\{\mathbf{y}_t\}$  is vector  $MA(\infty)$  with absolutely summable coefficients and  $\{\boldsymbol{\varepsilon}_t\}$  is vector independent white noise, then

$$\sqrt{n}(\bar{\mathbf{y}} - \boldsymbol{\mu}) \xrightarrow{d} N\left(\mathbf{0}, \sum_{j=-\infty}^{\infty} \boldsymbol{\Gamma}_j\right).$$

# HAC Covariance Matrix Estimator

- Recall that in Assumption 2.5, we assume that  $\{\mathbf{g}_t\}$  is an m.d.s. Since no serial correlation is allowed under this assumption, the asymptotic variance of  $\bar{\mathbf{g}}$  is simply  $\mathbf{S} = \mathbb{E}(\mathbf{g}_t \mathbf{g}_t')$ .
- If there is serial correlation in  $\{\mathbf{g}_t\}$ , the variance of the limiting distribution of  $\sqrt{n}\bar{\mathbf{g}}$  equals the long-run covariance matrix,

$$\mathbf{S} = \sum_{j=-\infty}^{\infty} \mathbf{\Gamma}_j = \mathbf{\Gamma}_0 + \sum_{j=1}^{\infty} (\mathbf{\Gamma}_j + \mathbf{\Gamma}_j'),$$

where  $\mathbf{\Gamma}_j$  is the  $j$ -th order autocovariance matrix

$$\mathbf{\Gamma}_j = \mathbb{E}(\mathbf{g}_t \mathbf{g}_{t-j}') \quad (j = 0, \pm 1, \pm 2, \dots).$$

# Estimating $\mathbf{S}$ When Autocovariances Vanish after Finite Lags

First, the consistent estimators of individual autocovariances are given by

$$\hat{\mathbf{\Gamma}}_j = \frac{1}{n} \sum_{t=j+1}^n \hat{\mathbf{g}}_t \hat{\mathbf{g}}_{t-j}' \quad (j = 0, 1, \dots, n-1)$$

where

$$\hat{\mathbf{g}}_t \equiv \mathbf{x}_t e_t.$$

Given the estimated autocovariances, and if we know that  $\mathbf{\Gamma}_j = \mathbf{0}$  for  $j > q$  where  $q$  is known and finite, then clearly  $\mathbf{S}$  can be consistently estimated by

$$\hat{\mathbf{S}} = \hat{\mathbf{\Gamma}}_0 + \sum_{j=1}^q (\hat{\mathbf{\Gamma}}_j + \hat{\mathbf{\Gamma}}_j') = \sum_{j=-q}^q \hat{\mathbf{\Gamma}}_j. \quad (14)$$

# Using Kernels to Estimate $S$

The **kernel-based** estimators can be expressed as a weighted average of estimated autocovariances:

$$\hat{S} = \sum_{j=-n+1}^{n-1} k\left(\frac{j}{q(n)}\right) \cdot \hat{\Gamma}_j.$$

Here, the function  $k(\cdot)$ , which gives weights for autocovariances, is called a kernel and  $q(n)$  is called the bandwidth. The estimator (14) is a special kernel-based estimator with  $q(n) = q$  and

$$k(x) = \begin{cases} 1 & \text{for } |x| \leq 1, \\ 0 & \text{for } |x| > 1. \end{cases}$$

This kernel is called the truncated kernel. This truncated kernel-based  $\hat{S}$ , however, is not guaranteed to be positive semidefinite in finite samples.

[Newey & West, 1987] noted that the kernel-based estimator can be made non-negative definite in finite samples if the kernel is the Bartlett kernel:

$$k(x) = \begin{cases} 1 - |x| & \text{for } |x| \leq 1, \\ 0 & \text{for } |x| > 1. \end{cases}$$

The Bartlett kernel-based estimator of  $\mathbf{S}$  is called the **Newey-West estimator**.

Under certain regularity conditions on random sample, kernel function  $k(\cdot)$ , and bandwidth  $q(n)$ , we have

$$\hat{\mathbf{S}} \xrightarrow[p]{} \mathbf{S},$$

provided  $q(n) \rightarrow \infty, q(n)/n \rightarrow 0$ .

# Clustered Sampling

In microeconometrics, the observations  $\{(y_1, \mathbf{x}_1), \dots, (y_i, \mathbf{x}_i), \dots, (y_n, \mathbf{x}_n)\}$  are assumed to be independent and identically distributed. This assumption may be violated if individuals in the sample are connected in some way, for example if they are

- neighbors,
- members of the same village,
- classmates at a school,
- firms within a specific industry.

A currently popular approach which allows for mutual dependence is known as **clustered dependence**.

In clustering contexts it is convenient to double index the observations as  $(y_{ig}, \mathbf{x}_{ig})$  where

- $g = 1, \dots, G$  indexes the cluster
- $i = 1, \dots, n_g$  indexes the individual within the  $g^{th}$  cluster

The total number of observations is  $n = \sum_{g=1}^G n_g$ .

Let  $\mathbf{y}_g = (y_{1g}, \dots, y_{n_gg})'$  and  $\mathbf{X}_g = (\mathbf{x}_{1g}, \dots, \mathbf{x}_{n_gg})'$  denote the  $n_g \times 1$  vector of dependent variables and  $n_g \times k$  matrix of regressors for the  $g^{th}$  cluster. A linear regression model can be written for the individual observations as

$$y_{ig} = \mathbf{x}'_{ig}\boldsymbol{\beta} + e_{ig}$$

and using cluster notation as

$$\mathbf{y}_g = \mathbf{X}_g\boldsymbol{\beta} + \mathbf{e}_g$$

where  $\mathbf{e}_g = (e_{1g}, \dots, e_{n_gg})'$  is a  $n_g \times 1$  error vector.



The OLS estimator of  $\beta$  can be written as

$$\begin{aligned}\hat{\beta} &= \left( \sum_{g=1}^G \sum_{i=1}^{n_g} \mathbf{x}_{ig} \mathbf{x}_{ig}' \right)^{-1} \left( \sum_{g=1}^G \sum_{i=1}^{n_g} \mathbf{x}_{ig} y_{ig} \right) \\ &= \left( \sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}_g' \mathbf{y}_g \right) \\ &= (\mathbf{X}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{y}).\end{aligned}$$

The residuals are  $\hat{e}_{ig} = y_{ig} - \mathbf{x}_{ig}' \hat{\beta}$  in individual level notation and  $\hat{e}_g = \mathbf{y}_g - \mathbf{X}_g \hat{\beta}$  in cluster level notation.

The standard clustering assumption is that the clusters are known to the researcher and that the observations are independent across clusters:

### Assumption

The clusters  $(\mathbf{y}_g, \mathbf{X}_g)$  are mutually independent across clusters  $g$ .

The model is a linear regression under the assumption

$$\mathbb{E}[\mathbf{e}_g | \mathbf{X}_g] = \mathbf{0},$$

or

$$\mathbb{E}[e_{ig} | \mathbf{X}_g] = 0.$$

Note that

$$\hat{\beta} - \beta = \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{e}_g \right).$$

We have

$$\begin{aligned} \mathbb{E}[\hat{\beta} - \beta | \mathbf{X}] &= \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbb{E}[\mathbf{e}_g | \mathbf{X}] \right) \\ &= \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbb{E}[\mathbf{e}_g | \mathbf{X}_g] \right) \\ &= \mathbf{0}. \end{aligned}$$

Therefore,

$$\mathbb{E}[\hat{\beta} | \mathbf{X}] = \beta.$$

Let

$$\Sigma_g = \mathbb{E} [e_g e_g' | \mathbf{X}_g]$$

denote the  $n_g \times n_g$  conditional covariance matrix of the errors within the  $g^{th}$  cluster. Since the observations are independent across clusters,

$$\begin{aligned} \text{Var} \left[ \left( \sum_{g=1}^G \mathbf{X}_g' e_g \right) | \mathbf{X} \right] &= \sum_{g=1}^G \text{Var} [\mathbf{X}_g' e_g | \mathbf{X}_g] \\ &= \sum_{g=1}^G \mathbf{X}_g' \mathbb{E} [e_g e_g' | \mathbf{X}_g] \mathbf{X}_g \\ &= \sum_{g=1}^G \mathbf{X}_g' \Sigma_g \mathbf{X}_g \equiv \Omega_n. \end{aligned}$$

It follows that

$$\mathbf{V}_{\hat{\beta}} = \text{Var}[\hat{\beta} | \mathbf{X}] = (\mathbf{X}' \mathbf{X})^{-1} \Omega_n (\mathbf{X}' \mathbf{X})^{-1}. \quad (15)$$

(15) differs from the formula in the independent case due to the correlation between observations within clusters. The magnitude of the difference depends on

- the degree of correlation between observations within clusters
- the number of observations within clusters

Consider a scenario in which  $n_g = N$ ,  $\mathbb{E}[e_{ig}^2 | \mathbf{X}_g] = \sigma^2$ ,  $\mathbb{E}[e_{ig}e_{\ell g} | \mathbf{X}_g] = \sigma^2\rho$  for  $i \neq \ell$ , and the regressors  $\mathbf{x}_{ig}$  do not vary within a cluster. In this case,

$$\mathbf{V}_{\hat{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \sigma^2(1 + \rho(N - 1)).$$

If  $N = 48$  and  $\rho = 0.25$ , the correct standard errors are a multiple of about **three times** the conventional formula.

$\Omega_n$  can be estimated by

$$\begin{aligned}
 \hat{\Omega}_n &= \sum_{g=1}^G \mathbf{X}_g' \hat{\mathbf{e}}_g \hat{\mathbf{e}}_g' \mathbf{X}_g \\
 &= \sum_{g=1}^G \sum_{i=1}^{n_g} \sum_{\ell=1}^{n_g} \mathbf{x}_{ig} \mathbf{x}_{\ell g}' \hat{e}_{ig} \hat{e}_{\ell g} \\
 &= \sum_{g=1}^G \left( \sum_{i=1}^{n_g} \mathbf{x}_{ig} \hat{e}_{ig} \right) \left( \sum_{\ell=1}^{n_g} \mathbf{x}_{\ell g} \hat{e}_{\ell g} \right)'.
 \end{aligned}$$

A natural cluster covariance matrix estimator takes the form

$$\hat{\mathbf{V}}_{\hat{\beta}} = a_n (\mathbf{X}'\mathbf{X})^{-1} \hat{\Omega}_n (\mathbf{X}'\mathbf{X})^{-1},$$

where  $a_n$  is a possible finite-sample adjustment. The Stata cluster command uses

$$a_n = \left( \frac{n-1}{n-k} \right) \left( \frac{G}{G-1} \right).$$

# References



Fumio Hayashi. (2000)

*Econometrics, Chapter 1, 2, 4 & 6.*

Princeton University Press, 2000.



Halbert White. (1980)

*A Heteroskedasticity-Consistent Covariance Matrix Estimator and a Direct Test for Heteroskedasticity.*

*Econometrica* **48**: 817–838.



Whitney K. Newey & Kenneth D. West (1987)

*A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix.*

*Econometrica* **55**: 703-708.